

Towards performance portability through locality-awareness for applications using one-sided communication primitives

Huan Zhou

High Performance Computing Center Stuttgart (HLRS)
University of Stuttgart
Germany
zhou@hlrs.de

José Gracia

High Performance Computing Center Stuttgart (HLRS)
University of Stuttgart
Germany
gracia@hlrs.de

Abstract—MPI is the most widely used data transfer and communication model in High Performance Computing. The latest version of the standard, MPI-3, allows skilled programmers to exploit all hardware capabilities of the latest and future supercomputing systems. The revised asynchronous remote-memory-access model in combination with the shared-memory window extension, in particular, allow writing code that hides communication latencies and optimizes communication paths according to the locality of data origin and destination. The latter is particularly important for today's multi- and many-core systems. However, writing such efficient code is highly complex and error-prone. In this paper we evaluate a recent remote-memory-access model, namely DART-MPI. This model claims to hide the aforementioned complexities from the programmer, but deliver locality-aware remote-memory-access semantics which outperforms MPI-3 one-sided communication primitives on multi-core systems. Conceptually, the DART-MPI interface is simple; at the same time it takes care of the complexities of the underlying MPI-3 and system topology. This makes DART-MPI an interesting candidate for porting legacy applications. We evaluate these claims using a realistic scientific application, specifically a finite-difference stencil code which solves the heat diffusion equation, on a large-scale Cray XC40 installation.

Keywords—DART-MPI; one-sided; application porting; data-locality

I. INTRODUCTION

Numerical simulation using distributed and parallel applications, e.g., studying aerodynamics, fluid dynamics, molecular dynamics, weather forecasting and so on, provides researchers with insight into complex natural phenomena. Various parallel programming models are developed for implementing efficient and scalable distributed simulations. The Message Passing Interface (MPI, [1]) remains as the dominant communication model as a result of its high performance, portability and standardization on cutting-edge computing systems.

Multi- or many-core processors are commonly deployed in today's computation clusters due to it fuels an explosive growth of processing capability. According to the new TOP500 Supercomputer report [2], the most powerful supercomputer as of June 2016 – Sunway – exemplifies this fact with core counts exceeding 10,000,000. However, taking into account additional

resource sharing, communication and synchronization within a node of multi- and many-core clusters poses new challenge for design of parallel applications. To get the optimal intra-node performance, researchers need to design their parallel programs to be aware of data-locality and exploit all available communication paths. For instance, data exchange within a node should bypass the MPI communication primitives and use direct memory load/store operations.

The MPI standard has been evolving to keep up with the scalable and productive computation hardware. MPI-2 incorporates the one-sided interfaces to avoid the matching affairs inherent in the two-sided operation. Moreover, in MPI-3 remote-memory-access (RMA) is extended to support shared memory windows. The results provided in earlier works [3], [4] have proved that the MPI-integrated shared memory window is a promising alternative since it allows co-existence of RMA operations and direct load/store accesses. However, note that MPI does not intuitively expose data locality to user applications. Thus, the programmer needs to have profound knowledge of the details in MPI to be able to write efficient, data-locality aware applications based on MPI. In addition, the MPI interfaces for RMA and shared-memory extensions are very complex, making the development error-prone and time-consuming.

DART-MPI [5] is an MPI-3 based implementation of the runtime system for the hierarchical PGAS-like C++ programming model DASH [6]. However, in this paper DART-MPI is treated as a programming model in its own right. It provides a number of communication primitives including one-sided and collective operations through a plain C-based interface. One fundamental goal of DART-MPI is to support data-locality aware communication schemes while keeping the user away from MPI complexities. In particular, DART-MPI internally takes responsibility over all those details related to the correctness and performance mentioned above. Therefore, DART-MPI enables portable and efficient parallel programming for today's multi- and many-core clusters. Besides, the DART RMA communication performance penalties for wrapping the underlying MPI-3 and system topology are shown

in the dissertation [7] and they become negligible when the larger messages are transferred.

We briefly discuss some background used in this paper in section II, explain how DART-MPI hides detail and complexities of MPI-3 RMA from the user and demonstrate the terseness and clarity of DART-MPI code in comparison to locality aware MPI-3 code in section III. Finally, in section IV we demonstrate the performance of DART-MPI using a heat diffusion problem, which is a representation for a wide class of numerical simulation schemes. In particular, we present a comprehensive picture of the comparative results between the DART-MPI and MPI-3 RMA implementations of this application on large-scale, supercomputers such as Cray XC40 systems.

II. BACKGROUND

A. DART-MPI

DART [5] defines common concepts, terminology and abstracts from the underlying communication substrate and hardware. DART establishes a partitioned global address space and provides functions to handle memory efficiently, such as memory allocation and data movement. In addition, DART also provides functions for initialization, synchronization and management of teams, which are similar to the MPI communicators. In a DART program an individual participate is called unit. DART provides the functions of *dart_init* and *dart_exit* for initialization and finalization. Importantly, providing and working with a global memory (collective or non-collective) is the focus of DART. The global address space is a virtual abstraction, with each unit contributing a part of its local memory. Remote data items are addressed by global pointers provided by DART. The complete DART specification is available on-line at https://dl.dropboxusercontent.com/u/408013/dart_spec_v2.1.html.

When possible (i.e., intra-node), the *shmem-win* (shared memory window object) described in the paper [8] are always used by DART-MPI blocking RMA operations to perform load or store instructions. Additionally, the DART-MPI intra-node non-blocking RMA operations are designed as the MPI RMA operations on *shmem-win*, whilst all the DART-MPI inter-node RMA operations turn to the MPI RMA operations on *d-win* (MPI dynamically-created window object) or *win* (MPI-created window object). Hence we can observe that DART-MPI spontaneously has data locality in mind when performing the RMA operations. DART blocking and non-blocking RMA operations are represented as *dart_get/put_blocking* and *dart_get/put*, respectively.

B. Heat conduction

The heat conduction [9] is a mode of heat transfer owing to molecular activity and occurs in any material (e.g., solids, fluids and gases). We used the application source code parallelized with MPI at [10]. This application simulates the phenomenon of 3D heat conduction in solids with temperature-dependent thermal diffusivity. The Partial Differential Equation (PDE) for *unsteady* 3D heat conduction is obtained in

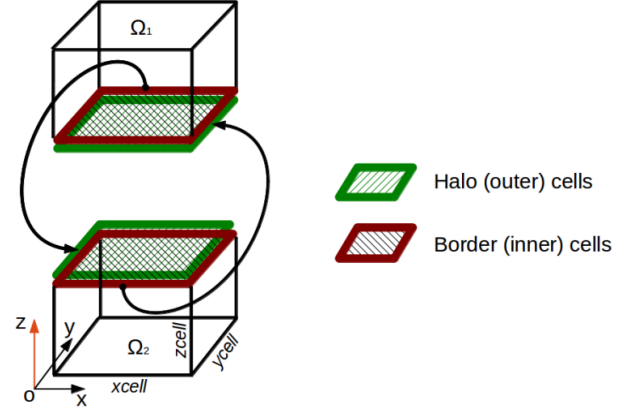


Fig. 1. Halo cells exchange in 3D grid. Exchange of the matrix data (a total of $x_{cell} \times y_{cell}$ grid cells) happen on Oxy plane between neighboring processes in the direction of z-axis, where sub-domain Ω_2 is the Down neighbor of sub-domain Ω_1 .

the Cartesian coordinate system. Here, *unsteady* means the temperatures may change during the process of heat conduction. This application solves the PDE over a 3D grid by using the Finite Difference Method (FDM). Boundary conditions are constant temperatures at the edges of the 3D grid.

Moreover, parallelization is done based on the checkerboard domain decomposition. After the decomposition each process has six neighbors located in the direction of *East-West* (x-axis), *North-South* (y-axis) and *Up-Down* (z-axis) according to the Cartesian process topology. The exceptional case is the processes owning the edge cells will not have all of the six neighbors. Each process sends each face of its sub-domain (a 3D sub-grid), i.e., matrix, to its corresponding neighbor for the computation.

Each process exchanges the border data with its corresponding neighbors and the halo cells are reserved to receive the exchanged data. The original 3D heat conduction application fills the halo cells in each iteration using the point-to-point communication operation (i.e., *MPI_Sendrecv*). Each cell of the grid is a 8-byte double-precision floating-point number. The abort-criterion convergence is achieved by invoking the collective operation, such as all-to-all reduce. Figure 1, takes the direction of z-axis for example, to explain the way of boundary data exchange between two neighboring processes.

III. METHOD OF WRAPPING MPI COMPLEXITY INSIDE DART INTERFACES

Figure 2 shows the MPI code that is used to consciously send a message from the origin to another random target process with the data locality in mind. It is clearly observed that in MPI code, the window creation/destroy operations as well as the access epoch start/end operations should be dealt with explicitly and carefully. In detail, the MPI programmer first needs to allocate a shared-memory window along with a region of memory with specified size and then creates another RMA window spanning the same region of memory. Besides the creation of these two nested windows, a hash table (named

```

MPI_Init (...);
MPI_Comm_size (MPI_COMM_WORLD, &size);
MPI_Comm_split_type
(MPI_COMM_WORLD,..., &sharedmem_comm);
MPI_Win_create_dynamic (MPI_COMM_WORLD, &d-win);
MPI_Win_lock_all (d-win);
if (sharedmem_comm!=MPI_COMM_NULL){
    MPI_Win_allocate_shared
    (nbytes,..., &membase, &sharedmem_win);
    MPI_Win_lock_all (sharedmem_win);}
MPI_Win_attach (d-win, membbase, nbytes);
MPI_Get_address (membbase, &disp);
disp_s = (MPI_Aint*)malloc (size *
    (sizeof(MPI_Aint)));
MPI_Allgather (&disp, 1, MPI_AINT, disp_s,...,
    MPI_COMM_WORLD);
dest = randomdst_generator ();
/* The array is_shmem indicates the position of
the given target process with respect to the
origin */
if ((j=is_shmem[dest])>=0){//on-node communication
/* The j is the relative target rank in the
corresponding sharedmem_comm */
MPI_Win_shared_query
(sharedmem_win, j,..., &baseptr);
memcpy (baseptr, srcptr, nbytes);}
else{//across-node communication
MPI_Put (srcptr, nbytes, dest, disp_s[dest],...,
    d-win);
MPI_Win_flush (dest, d-win);}
MPI_Win_detach (d-win, membbase);
if (sharedmem_comm!=MPI_COMM_NULL){
    MPI_Win_unlock_all (sharedmem_win);}
MPI_Win_unlock_all (d-win);
MPI_Win_free (&sharedmem_win);
MPI_Win_free (&d-win);
free (disp_s);
MPI_Finalize ();

```

Fig. 2. An example for MPI code. This example shows the MPI code for transferring data from an origin process to another random target process with data locality in mind.

```

dart_init ();
dart_team_memalloc_aligned
(DART_TEAM_ALL, nbytes, &gptr);
dest = randomdst_generator ();
dart_gptr_setunit (&gptr, dest);
dart_put_blocking (gptr, srcptr, nbytes);
dart_team_memfree (DART_TEAM_ALL, gptr);
dart_exit ();

```

Fig. 3. An example for DART code. This example shows the DART code for transferring data from an origin unit to another random target unit with data locality in mind.

as `is_shmem`) is entailed to indicate the data locality. Finally, different communication paths are taken according to the data locality information (on-node or across-node). Regarding the on-node communications, direct load or write accesses are allowed. Otherwise, the MPI RMA communication operations are employed. Figure 3 shows the corresponding data-locality aware code of DART. Revisiting Fig. 2, the red code lines are required for handling the intra-node data transfers. In comparison to the MPI code, the DART code is obviously more concise and easier-to-read. This is due to that DART-MPI has internally taken over the responsibility for the locality-

awareness and manipulation of access epoch by hiding the complexity from the user.

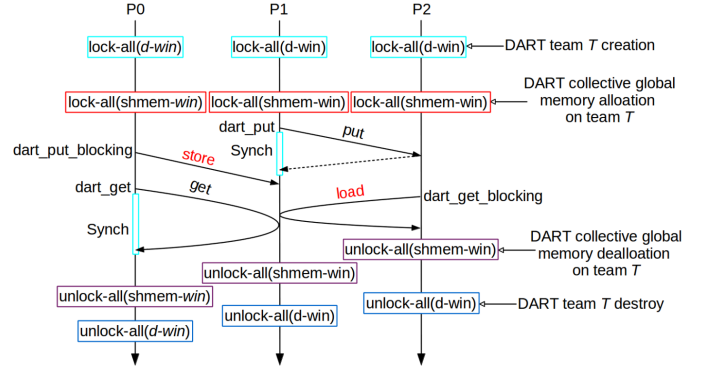


Fig. 4. Example of DART-MPI RMA operations within two passive epochs. Here, the two passive epochs are started by a call to `MPI_Win_lock_all` and ended by a call to `MPI_Win_unlock_all`. The `Synch` means the operation of `dart_wait`.

The MPI RMA communication operations coupled with explicit passive synchronization operations (with shared lock) form a theoretic foundation for building the DART-MPI RMA model [5]. In fact, each MPI RMA communication call must proceed within an access epoch for this window at a process to start and complete the issued RMA communications. Superficially, the RMA semantics in MPI and DART both require the independent management of synchronization and communication. However, the concept of access epoch is not relevant in the DART RMA semantics according to the DART specification. Therefore, all primitives pertaining to MPI access epoch need to be hidden from DART programmer. Importantly, this focus creates the illusion that all DART global memory regions are protected and exposed once allocated and also the expected processes are allowed to access those global memory regions directly.

MPI-3 supports a global lock mechanism providing the `lock_all/unlock_all` routines (i.e., `MPI_Win_lock_all/MPI_Win_unlock_all`) to allow locking/unlocking multiple targets simultaneously. I.e., the pair of `MPI_Win_lock_all` and `MPI_Win_unlock_all` enables programmer to lock/unlock all processes in a certain window with a shared lock. This mechanism avoids the proneness to error and repeated open or closure operations on the access epoch. Thus, it is applied to build up DART-MPI RMA model by leveraging the global lock/unlock routines. Given the DART collective and non-collective global memory managements differ in concept and design, we will explain how to safely expose these two types of global memory independently based on the global lock mechanism.

A global `win` and a global `shmem-win` are created once a DART program is initiated [8]. Therefore, each process/unit needs to add two `MPI_Win_lock_all` calls for the global `win` and `shmem-win` to start two protected shared access epochs to all other units. Prior to freeing up the global `win` and `shmem-win` (done inside `dart_exit`), each unit issues two matched

MPI_Win_unlock_all calls to end the above two access epochs.

Unlike non-collective global memory allocations, collective global memory allocations always involve all units in the given team. Thus, there are two scenarios where the creation of *d-win* [8] takes place:

- 1) When a DART program is launched, a *d-win* for *DART_TEAM_ALL* is created.
- 2) When a new team (e.g., team *T*) is created, a *d-win* for team *T* is created.

In each of the above scenarios, a shared lock all epoch for the *d-win* is started. Additionally, we associate a *shmem-win* with a region of allocated collective global memory. This necessitates the call to *MPI_Win_lock_all* for the *shmem-win*.

Likewise, we can complete the lock all epoch by a call to *MPI_Win_unlock_all* when the DART program is finalized or team *T* is destroyed. A call to *MPI_Win_unlock_all* is also entailed to complete a shared RMA access epoch along with the deallocation of the collective global memory region.

Remote completion of communications without ending the access epoch can be achieved with the *flush* routines (i.e., *MPI_Win_flush* and its all-, any- and local-variants). An *MPI_Win_flush* specifies a certain target and ensures that all previous operations to the target are locally and remotely finished.

The DART synchronization routines (*dart_wait* and *dart_waitall*) are expected to ensure the remote and local completion of DART non-blocking RMA communications. Therefore, the explicit MPI bulk synchronization using *flush* is integrated into *dart_wait*. Likewise, the *dart_waitall* performs a series of related calls to *MPI_Win_flush*. Figure 4 thoroughly shows the exemplary DART RMA events happening on the collective global memory region across team *T* with MPI global lock mechanism. Here we assume that the team *T* consists of three processes locating within one node.

IV. EXPERIMENTAL EVALUATION

In this section we present a comprehensive picture of the comparative results between DART-MPI RMA and MPI RMA through a numerical simulation of 3D heat conduction problem. In this experiment, we stop the calculation after 5000 iterations and collect and analyze the time results of the computation (update the inner grid cells) and halo cells exchange for different grid sizes and participating processes. We plot the average data results of 25 runs with small execution time variation reported.

A. Experimental testbed

The experimental environment was the Cray XC40 system [11]. Cray XC40 system is equipped with 7,712 compute nodes made up of dual twelve-core Intel Haswell E5-2680v3 processor (one processor per socket), which has exclusive 256 KB L2 unified cache for each core. Therefore, each compute node has 24 cores running at 2.5GHZ with 128 GB of DDR4 (Double Data Rate) RAM (Random Access Memory). The different compute nodes are interconnected via a Cray Aries network using Dragonfly topology. The Cray XC40

```

/* Communicate matrix data on Oxz plane */
if (South neighbor exists)
  for (i = 0; i < xcell; i++)
    get (zcell, South); // Get zcell consecutive
                        // grid cells from South neighbor
barrier;
if (North neighbor exists)
  for (i = 0; i < xcell; i++)
    get (zcell, North);
barrier;
/* Communicate matrix data on Oyz plane */
if (West neighbor exists)
  for (i = 0; i < ycell; i++)
    get (zcell, West); // Get zcell consecutive grid
                        // cells from West neighbor
barrier;
if (East neighbor exists)
  for (i = 0; i < ycell; i++)
    get (zcell, East);
barrier;
/* Communicate matrix data on Oxy plane */
if (Up neighbor exists)
  for (i = 0; i < xcell; i++)
    for (j = 0; j < ycell; j++)
      get (1, Up); // Get one grid cell from Up
                  // neighbor
barrier;
if (Down neighbor exists)
  for (i = 0; i < xcell; i++)
    for (j = 0; j < ycell; j++)
      get (1, Down);
barrier;

```

Fig. 5. Pseudo-code for the halo cells exchange with get and barrier operations.

system features a hierarchical network topology. Detailedly, the Cray XC Rank-1 network is used for communications happening across different compute blades within a chassis over backplane. Each chassis consists of 16 compute blades. Communications happening across different chassis within a 2-cabinet group over copper cables are supported with the Cray XC Rank-2 network. Further, the Cray XC Rank-3 network is employed for communications happening between the groups over the optical cables. All the experiments in this paper are based on the Cray Programming Environment 5.2.56. Moreover, the tasks are assigned to cores in SMP style [12], which means one node needs to be filled before going to next.

B. DART-MPI and MPI RMA port of 3D heat conduction application

For our test, we first implement the above mentioned 3D heat conduction algorithm based on MPI one-sided interfaces using passive target mode. We then port it to DART one-sided directives, where just several (~ 14) modified code lines are required. The porting procedure is thus straightforward and economic. The halo cells exchange is achieved by using blocking get operation. We can deduce that six blocking get operations will be invoked for each process. Particularly, *MPI_Rget* is invoked first and then followed by *MPI_Wait* in MPI implementation and *dart_get_blocking* is used in DART implementation. Note that, within each calculation iteration, after one round of halo cells exchange with certain neighbor an

explicit process synchronization, such as barrier (*MPI_Barrier* in MPI and *dart_barrier* in DART-MPI), is naively inserted in our code due to an implicit synchronization is ideally supported by point-to-point communication model. Therefore, the halo cells exchange time measured below includes the overhead brought by process synchronization. Regarding the process synchronization, we only measure the overhead introduced by the barrier operations instead of the time to wait between processes. Figure 5 concisely shows the critical part of the halo cells exchange code within each calculation iteration.

C. Performance results

Figure 6 gives the quantitative results showing the performance of DART-MPI RMA and MPI RMA on Cray XC40 system over the number of cores (from 16 to 4096). A weak scaling evaluation, where the number of grid cells per core is fixed and is applied with grid sizes varying from $(64 \times 64 \times 128)$ to $(256 \times 2048 \times 256)$. The domain size is chosen such, that overall execution is dominated by the time to exchange halo cells in order to increase the statistical significance of the latter. As expected from a weak scaling experiment, the total execution time increases only slowly due to larger relative communication cost. In fact, the actual computation time, i.e. the time to update the inner grid cells, stays nearly constant with the increasing core count and is the same for the MPI RMA and DART-MPI version within the observed statistical variance. However, the DART-MPI implementation can provide a relative speedup of $\sim 1.38\times$ over MPI on average in terms of the halo cells exchange time performance. Such speedup is mostly attributed to the enabling of direct load/store access within one node in the DART-MPI. Particularly, when only the on-node performance is considered, that is 16 cores, then the improvement can be obviously seen, cf. 0.76s for DART-MPI, and 1.85s for MPI. In addition, the consistent improvement seen is expected because most of the data exchanges are still local to a node in this evaluation. Such fact highlights and visualizes the advantage of the memory sharing mechanism in intra-node case employed by DART-MPI implementation. Furthermore, we can observe that the performance speedup gets decreased as the number of cores is increased, which is not surprising given the number of across-node data exchanges is accordingly increased and DART RMA uses the same communication infrastructure as MPI RMA internally in the inter-node case.

On the other hand, shown in Fig. 6, the halo cells exchange time gets sustained growth over the number of cores. Specifically, a sudden rise occurs when the communications go beyond one node and start crossing nodes (32 cores). The fact, that is, several inter-node data transfers start to be involved accordingly, creates big difference in performance. Besides that, the process synchronization may also be a part of the growth of the halo cells exchange time, as hinted above. Therefore, Figure 7 breaks down the halo cells exchange time in terms of pure data exchange and process synchronization, which can help us understand how the two components affect

the amount of time the halo cells exchange takes to run. Observing this figure, it is immediately clear that both the process synchronization and pure data exchange overhead increase gradually over the number of cores. The breakdown information tells us that the pure data exchange component plays an important role in the halo cells exchange time rise especially when the number of cores reaches up to 1024. In this experiment, inter-group communications through Cray XC Rank-3 optical network take place when 1024 cores are launched, which would greatly degrade data exchange performance, in comparison to the inter-blade (intra-chassis) communications through Cray XC Rank-1 backplane network and inter-chassis (intra-group) communications through Cray XC Rank-2 copper network. This is the reason for the sudden rise in the pure data exchange time when ranging the number of cores from 512 to 1024.

Here, we would add that the Fig. 7 sufficiently illustrates that the data-locality aware RMA operations in the DART-MPI result in the improvement in the halo cells exchange time.

V. RELATED WORK

A heuristic study [13] conducted on Cray XE6 shows that one-sided MPI-2 has relatively inferior scaling behavior compared to UPC and Cray SHMEM RMA. This is due to that random allocated memory can be specified to be remotely accessible in MPI-2. Unlike MPI-2, the memory for RMA operation is somehow customized for Cray SHMEM (symmetric memory) and UPC (shared memory). Further, MPI-3 RMA enables direct load/store accesses by supporting a shared-memory window [3], [4], which is employed in DART-MPI to enhance the intra-node RMA communication performance.

Additionally, MPI researchers conduct long-term optimization works on the MPI RMA [14]–[17] and collective operations [18]–[21]. Those optimizations are also liable to be applied onto DART-MPI and then benefit the performance of applications.

VI. CONCLUSION

We depicted the approach of leveraging the existing MPI RMA synchronization interfaces (global lock mechanism) for the related DART interfaces. The DART code of data locality-awareness is more concise and easier-to-read than MPI code. This is due to that DART-MPI has internally taken over the responsibility for the locality-awareness and manipulation of epoch access. We then described the rewriting structures for porting a 3D heat conduction application onto one-sided MPI and then onto DART-MPI. The DART-MPI port is proved to be efficient for different number of cores and grid sizes on Cray XC40 system and speeds up the communication performance (halo cells exchange) by 27% on average over MPI port in this application level evaluation. Such performance improvement is induced by the fact that one-sided DART is implemented using the feature of MPI-3 integrated shared memory window in the intra-node case while retaining the MPI inter-node RMA performance. Therefore, DART-MPI has basically reached

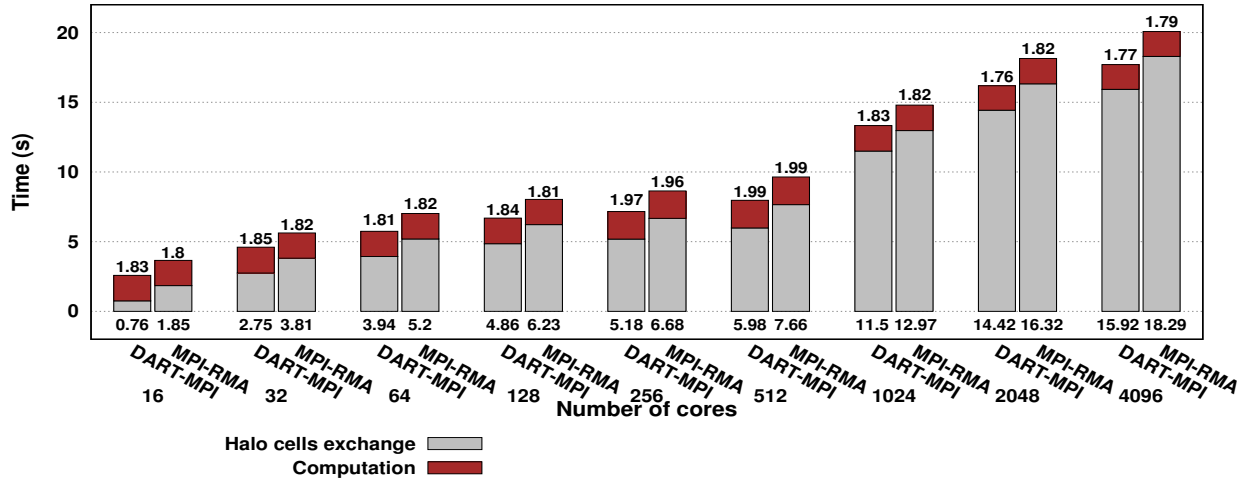


Fig. 6. Weak scaling for one-sided DART and MPI. The grey bar signifies the halo cells exchange time, the value of which is written on the bottom. The brown bar signifies the computation overhead, the value of which is written on the top.

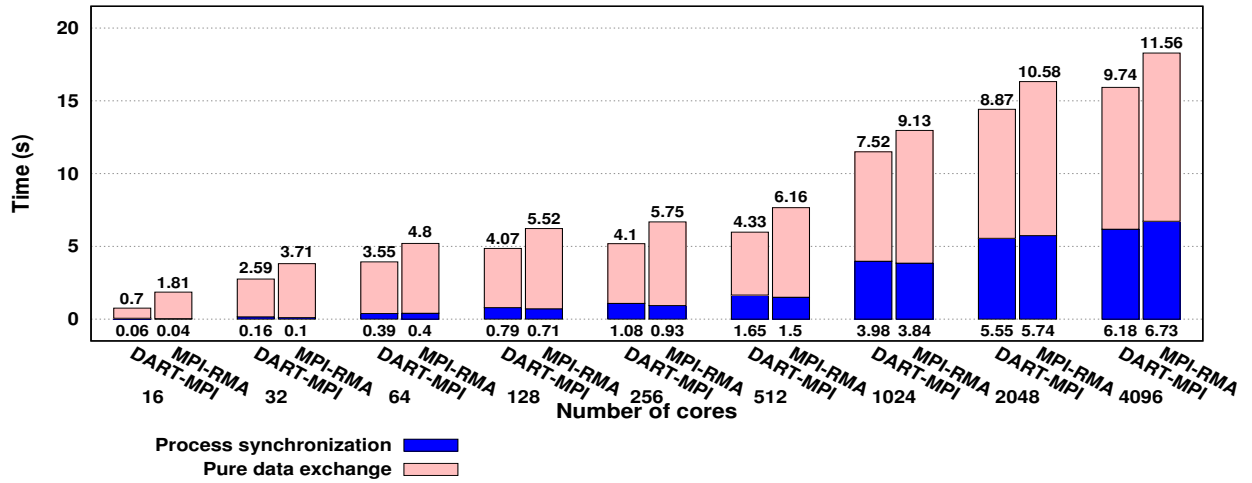


Fig. 7. Breakdown of halo cells exchange time for the weak scaling problem. The blue bar signifies the overall process synchronization overhead, the value of which is written on the bottom. The pink bar signifies the pure data exchange overhead, the value of which is written on the top. The sum equals to the halo cells exchange overhead that we measured in this benchmark.

the goal of making itself easier for programmers to write efficient applications with the data locality in mind. Further, the DART-MPI applications can be executed on different computers with hierarchical memory architecture (like Cray XC40 system) without any code change. Meantime, the performance advantages brought by the factor of locality awareness will also be retained since the underlying MPI communication conduit is highly supported and tuned by a wide range of hardware systems.

ACKNOWLEDGMENT

This work has been supported by the European Community through the project Polca (FP7 programme under grant agreement number 610686) and Mont Blanc 3 (H2020 programme under grant agreement number 671697). We gratefully acknowledge funding by the German Research Foundation (DFG) through the German Priority Programme 1648 Software for Exascale Computing (SPPEXA) and the project DASH.

REFERENCES

- [1] MPI Forum, “MPI: A Message-Passing Interface Standard. Version 3.0,”

- Tech. Rep., Sep. 2012, available at: <http://www.mpi-forum.org>.
- [2] "TOP500 - List Statistics - Jun. 2016," <http://www.top500.org/statistics/list/>, accessed: Jun. 2016.
 - [3] T. Hoefer, J. Dinan, D. Buntinas, P. Balaji, B. W. Barrett, R. Brightwell, W. Gropp, V. Kale, and R. Thakur, "Leveraging MPI's One-Sided Communication Interface for Shared-Memory Programming," in *EuroMPI*, ser. Lecture Notes in Computer Science, J. L. Träff, S. Benkner, and J. J. Dongarra, Eds., vol. 7490. Springer, 2012, pp. 132–141. [Online]. Available: <http://dx.doi.org/10.1007/978-3-642-33518-1>
 - [4] J. R. Hammond, S. Ghosh, and B. M. Chapman, "Implementing OpenSHMEM Using MPI-3 One-Sided Communication," in *OpenSHMEM*, ser. Lecture Notes in Computer Science, S. W. Poole, O. R. Hernandez, and P. Shams, Eds., vol. 8356. Springer, 2014, pp. 44–58.
 - [5] H. Zhou, Y. Mhedheb, K. Idrees, C. W. Glass, J. Gracia, and K. Furlinger, "DART-MPI: An MPI-based Implementation of a PGAS Runtime System," in *PGAS*, A. D. Malony and J. R. Hammond, Eds. ACM, 2014, pp. 3:1–3:11. [Online]. Available: <http://dl.acm.org/citation.cfm?id=2676870>
 - [6] K. Furlinger, C. W. Glass, J. Gracia, A. Knüpfer, J. Tao, D. Hünich, K. Idrees, M. Maiterth, Y. Mhedheb, and H. Zhou, "DASH: Data Structures and Algorithms with Support for Hierarchical Locality," in *Euro-Par 2014: Parallel Processing Workshops - Euro-Par 2014 International Workshops, Porto, Portugal, August 25-26, 2014, Revised Selected Papers, Part II*, 2014, pp. 542–552.
 - [7] H. Zhou, "Communication Methods for Hierarchical Global Address Space Models in HPC (Unpublished)," Ph.D. dissertation, University of Stuttgart, Germany, 2016.
 - [8] H. Zhou, K. Idrees, and J. Gracia, "Leveraging MPI-3 Shared-Memory Extensions for Efficient PGAS Runtime Systems," in *Euro-Par*, ser. Lecture Notes in Computer Science, J. L. Träff, S. Hunold, and F. Versaci, Eds., vol. 9233. Springer, 2015, pp. 373–384. [Online]. Available: <http://dx.doi.org/10.1007/978-3-662-48096-0>
 - [9] M. Williams, "What is heat conduction?" <http://phys.org/news/2014-12-what-is-heat-conduction.html>, Dec. 2014.
 - [10] "MPI Parallelization for numerically solving the 3D Heat equation," https://dournac.org/info/parallel_heat3d, accessed: Apr. 2016.
 - [11] "Cray XC40," https://wickie.hlr.de/platforms/index.php/CRAY_XC40_Hardware_and_Architecture, accessed: Jul. 2016.
 - [12] "Reordering MPI Ranks," <http://www.nersc.gov/users/computational-systems/retired-systems/hopper/performance-and-optimization/reordering-mpi-ranks/>, accessed: Apr. 2016.
 - [13] C. M. Maynard, "Comparing One-sided Communication with MPI, UPC and SHMEM," in *Proceedings of the Cray User Group (CUG)*, 2012.
 - [14] J. Liu, W. Jiang, P. Wyckoff, D. K. Panda, D. Ashton, D. Buntinas, W. D. Gropp, and B. R. Toonen, "Design and Implementation of MPICH2 over InfiniBand with RDMA Support," in *IPDPS*. IEEE Computer Society, 2004. [Online]. Available: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=9132>
 - [15] P. Lai, S. Sur, and D. K. Panda, "Designing Truly One-sided MPI-2 RMA Intra-node Communication on Multi-core Systems," *Computer Science - R&D*, vol. 25, no. 1-2, pp. 3–14, 2010.
 - [16] R. Gerstenberger, M. Besta, and T. Hoefer, "Enabling Highly-scalable Remote Memory Access Programming with MPI-3 One Sided," in *SC*, W. Gropp and S. Matsuoka, Eds. ACM, 2013, pp. 53:1–53:12. [Online]. Available: <http://doi.acm.org/10.1145/2503210>
 - [17] S. Potluri, H. Wang, V. Dhanraj, S. Sur, and D. K. Panda, "Optimizing MPI One Sided Communication on Multi-core InfiniBand Clusters Using Shared Memory Backed Windows," in *EuroMPI*, ser. Lecture Notes in Computer Science, Y. Cotronis, A. Danalis, D. S. Nikolopoulos, and J. Dongarra, Eds., vol. 6960. Springer, 2011, pp. 99–109. [Online]. Available: <http://dx.doi.org/10.1007/978-3-642-24449-0>
 - [18] R. Thakur, R. Rabenseifner, and W. Gropp, "Optimization of Collective Communication Operations in MPICH," *IJHPCA*, vol. 19, no. 1, pp. 49–66, 2005.
 - [19] J. Liu, A. R. Mamidala, and D. K. Panda, "Fast and Scalable MPI-Level Broadcast Using InfiniBand's Hardware Multicast Support," in *IPDPS*. IEEE Computer Society, 2004. [Online]. Available: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=9132>
 - [20] A. R. Mamidala, R. Kumar, D. De, and D. K. Panda, "MPI Collectives on Modern Multicore Clusters: Performance Optimizations and Communication Characteristics," in *CCGRID*. IEEE Computer Society, 2008, pp. 130–137. [Online]. Available: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=4534181>
 - [21] H. Zhou, V. Marjanovic, C. Niethammer, and J. Gracia, "A Bandwidth-saving Optimization for MPI Broadcast Collective Operation," in *Proceedings of the International Conference on Parallel Processing Workshops, ICPPW*, Sep. 2015.
 - [22] *18th International Parallel and Distributed Processing Symposium (IPDPS 2004), CD-ROM / Abstracts Proceedings, 26-30 April 2004, Santa Fe, New Mexico, USA*. IEEE Computer Society, 2004. [Online]. Available: <http://ieeexplore.ieee.org/xpl/mostRecentIssue.jsp?punumber=9132>